

# Probability Part of Sigmoid Neuron

**M** [medium.com/@manveetdn/notes-on-basic-probability-part-of-sigmoid-neuron-padhai-onefourthlabs-course-a-first-course-on-c6a4038b072e](https://medium.com/@manveetdn/notes-on-basic-probability-part-of-sigmoid-neuron-padhai-onefourthlabs-course-a-first-course-on-c6a4038b072e)

**Disclaimer:** This is notes on “*Probability Part of Sigmoid Neuron*” Lesson  
(PadhAI onefourthlabs course “A First Course on Deep Learning”)



Denoting Universal and individual sets.

1. For any event  $A$   $P(A) \geq 0$  always probability any event lies between 0 and 1.
2. If  $A_1, A_2, A_3, \dots, A_n$  are disjoint (dis joint means  $A_i \cap A_j = \emptyset \forall i \neq j$ )
3. **Union of all  $i$  events probability = Sum of all individual probabilities of  $i$  events.**
4. If  $\Omega$  is the universal set (set containing all sets)
5.  **$P(\Omega) = 1$**  (Probability of universal set)

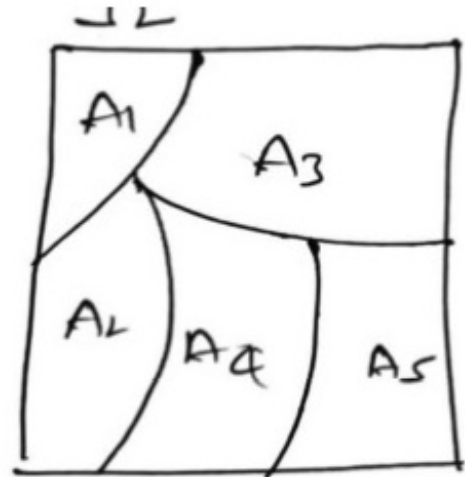
## Example :

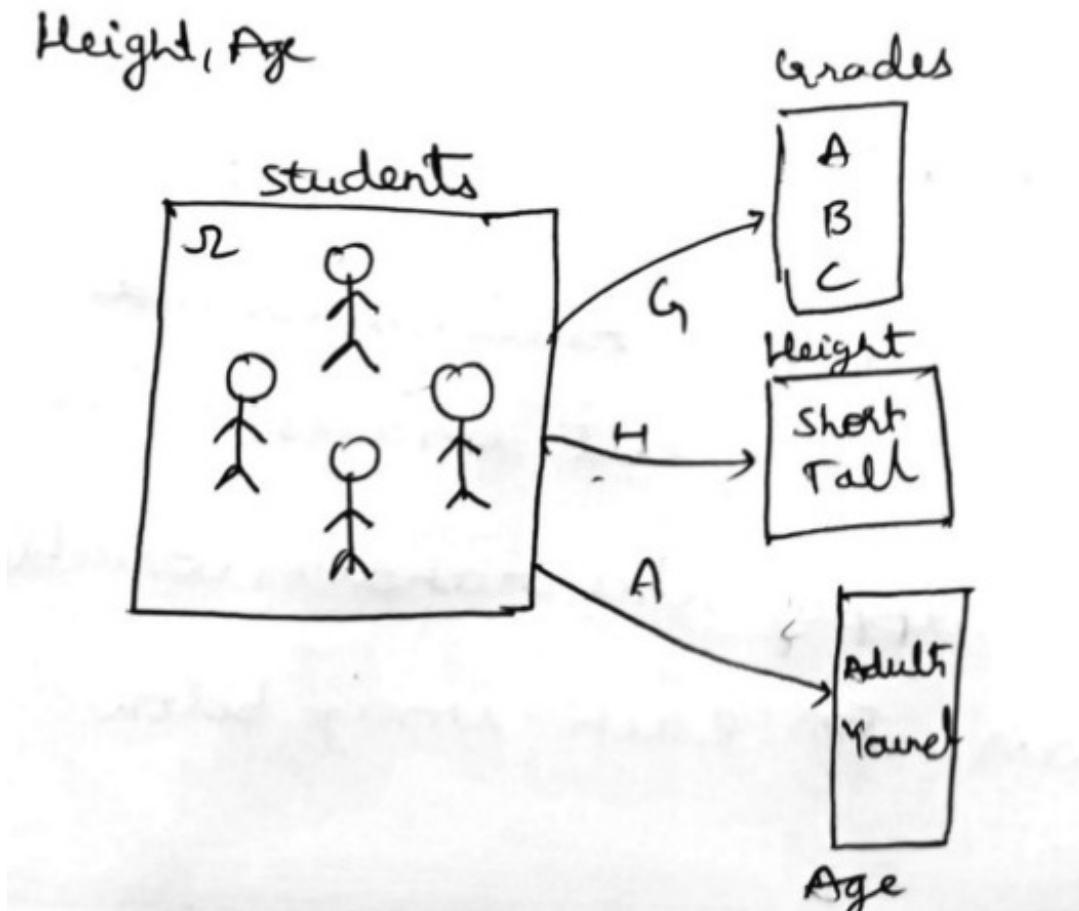
Suppose a student gets one  $S$  possible grades in a course : **A, B, C**

$P(\text{grade} = A), P(\text{grade} = B), P(\text{grade} = C)$

Actually in this case normally probabilities is calculated as (Numbers of 'A' Grades)/(Total number of students)

Actually coming to students we can have **multiple properties** associated with like **Grades, Height, Age.**





Each multiple-properties related to each student.

### Random Variable:

A random variable is a function which maps each outcomes  $\Omega$  to a value.

In the students example  $G$  (or  $f_{\text{grade}}$ ) maps each student in  $\Omega$  to a value: **A, B or C**.

The event  $\text{Grade} = A$  is a shorthand for the event  $\{\omega \in \Omega: f_{\text{grade}} = A\}$ .

Actually here we can take capstone project as example is that we have a heap of images and we need to decide whether it has text or no text

Therefore, We will initialise a folder having all images.

And if we apply this random variable called class on each image below.

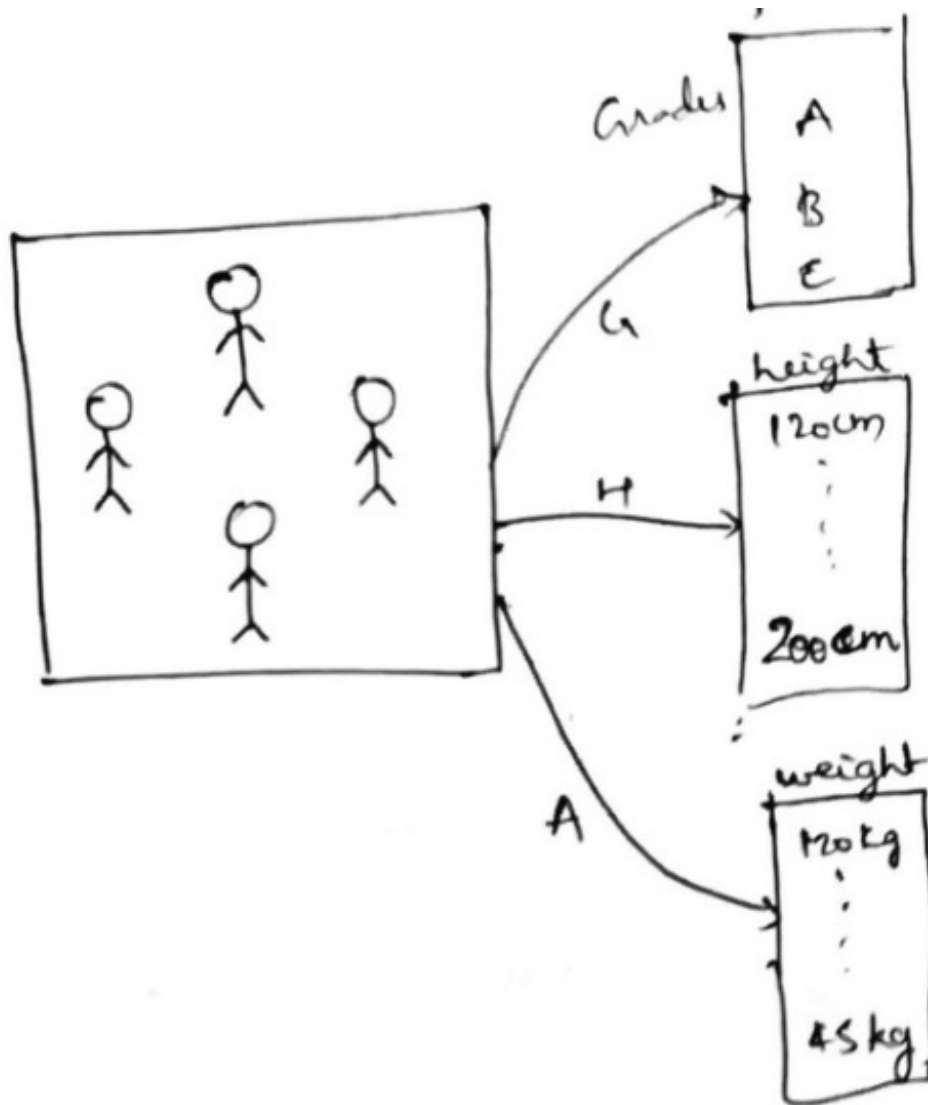
$f_{\text{class}}(\text{image}) = \begin{matrix} \text{Text} & 1 \\ \text{or} & \\ \text{No Text} & 0 \end{matrix}$   
 applying random variable on class image

What random variable variable says is that it takes all elements i.e., each and every element in the set and classify it as **Text or noText(1 or 0)**.

Now we can calculate **P(class = 0)** [this is short hand of above applying random variable on image and deciding whether there is text or not]

For example, we are given letter अ <sup>not.</sup>  
 here it comes in two case the letter is  
 a vowel or a consonant It contains text  
 or no text and also what is type of text  
 $P(\text{consonant} = \text{अ}) = ? \quad \cdot \quad P(\text{vowel} = \text{अ}) = ?$

Probability checking.



A random variable can either take continuous values for example the height and weight like **height ranges from 120 cm to 200 cm and weight ranges from 120 Kg to 45 Kg** or they may take **discrete values** like grades as it ranges on **A or B or C or D**.

In this we mainly focus on discrete random variable eg: **Grades(A to c), Rating(1 to 5), Vowels(a to u)**.

### **Distribution:(Marginal Distribution):**

Lets the random variable be A,B,C and the what is the distribution of grades G is given as below table which we call as distribution(Marginal Distribution).

Specify the marginal distribution over  $G$ ,  
means Specifying  $P(G=g) \forall g \in A, B, C$

We denote marginal distribution compactly  
by  $P(G)$ .

Let take bag full of balls R(Red) G(Green)  
B(Blue) colours.

Therefore, from this probability of ball being  
red is **(Number of red balls / Total number  
of balls.)**

Else probabilities of R, G, B are  
0.25, 0.4, 0.35

If you allow your friend to peep into the  
(bag) earn coating balls he will some how  
guess by estimating them and guess the  
probability when he is allowed to peep  
into it once and asked to estimate  
probabilities of each

Therefore, That is the predicted probability  
 $\hat{y}$ .

$G$	$P(G=g)$
A	0.1
B	0.2
C	0.7

R	G	B
0.25	0.4	0.35

$X$	$P(X=c)$
R	0.25
G	0.4
B	0.35

R G B

$y = [0.25 \ 0.4 \ 0.35] \rightarrow \text{true } y \text{ values}$

$\hat{y} = [0.3 \ 0.3 \ 0.4] \rightarrow \text{predicted } y \text{ values}$

True and Predicted values

We can directly say that he is work in guessing but we wanna say how wrong he was in  
predicting.

Then you can treat **y** and **yhat** calculate square error loss.

Certain event:

The probability of event if it is sure that it happens then it is called Sure Event and the probability of the event is 1.

A B C D  
[ 1 0 0 0 ]

Probability if A winning is certain event

If in match played by A,B,C,D then A has won.

A B C D  
[ 0.6 0.2 0.15 0.05 ]

If anyone comes suddenly you have other case where next match you can say I have been Observing the match from 1hr or so the **probability of A winning is 0.6 and that of B is 0.2 , C is 0.15 and D is 0.05 .**

Like this we can say in different ways the above are the two different ways of saying.

MUMBAI →  $\begin{matrix} \text{NoText} & \text{Text} \\ [ 0 & 1 ] \end{matrix}$  y

In the context of capstone project we will be given an image with text **The basic job is say that It contains Text or not** , the class has only two possibilities whether **text or NoText**.

What we will do at training time is that we will use sigmoid function we need output to **1 If it contains text and 0 if it doesn't contain text**. We will get the output zero or one like that

Assume 30x30 Image. Therefore, 900 inputs and therefore, 90 bias.

Therefore, if  $\hat{y} = 0.7$  is the thing which we got.

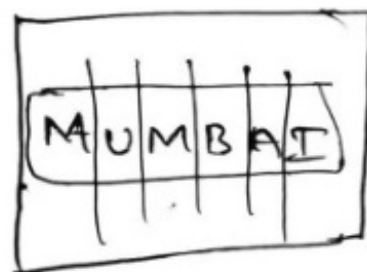
**Therefore,  $\hat{y} = 0.7$  is a probability distribution and it is saying the 0.7 probability it contains text and  $(1-0.7) = 0.3$  probability it doesn't contain text.**

If the model is perfect it **would have returned 1 for having Text and 0 for not having text.**

Now use square error loss or we will use **Cross Entropy loss which is more grounded for probability.**

This is all for binary classification. Whenever it comes to **multi class classification** is as below

The whole word is segmented into all single letters. After segmenting each one is checked with 26 alphabets in English and compared



∴ This is called 1-hot vector in which only one entry is 1 and everything else is zero

C	P(C=?)
a	0
b	0
...	...
m	1
...	...
x	0
y	0
z	0

$$\hat{y} = [\leftarrow 0000 \rightarrow 1 \rightarrow 000 \rightarrow] \mathbb{R}^{26}$$

This is true distribution and predicted distribution looks like

$$\hat{y} = [\leftarrow 0.01 \rightarrow 0.7 \rightarrow 0.3 \leftarrow]$$

∴ That is about the true & predicted distributions is probability

Now again use Square Error loss

$$\sum_{i=1}^{26} (y_i \hat{y}_i)^2$$

The probability distribution for the Looks like as below for each letter **let the letter be 'm'**. Then all **other letters probabilities will be zero and only 'm' letter probability will be one.**

This is a **1-hot vector** in which **only one entry is 1** and **everything else is zero.**

Like this Probability plays a major role in building a Deep Learning Model.

This is a small try, uploading the notes . I believe in **"Sharing knowledge is that best way of developing skills"**. Comments will be appreciated. Even small edits can be suggested.

Each Applause will be a great encouragement.

**Do follow my medium for more updates.....**